

DOI: 10.17148/IJIREEICE.2025.131033

# MACHINE LEARNING-BASED NETWORK INTRUSION DETECTION SYSTEM USING THE CSE-CIC-IDS2018 DATASET

# Sahithyaa Krishna Kumar<sup>1</sup>, R Sivani<sup>2</sup>, M Jaiaakash<sup>3</sup>, Dr Golda Dilip<sup>4</sup>

Student, Dept. of CSE, SRM Institute of Science and Technology, Chennai<sup>1</sup>

Student, Dept. of CSE, SRM Institute of Science and Technology, Chennai<sup>2</sup>

Student, Dept. of CSE, SRM Institute of Science and Technology, Chennai<sup>3</sup>

Guide, Professor Dept. of CSE, SRM Institute of Science and Technology, Chennai<sup>4</sup>

Abstract: Network Intrusion Detection Systems (NIDS) are vital defences against evolving and sophisticated cyber threats. Traditional security approaches frequently fail to detect novel, low-volume polymorphic attacks, necessitating the integration of adaptive machine learning (ML) models. This paper presents a high-performance, computationally efficient ML-based NIDS utilizing the contemporary CSE-CIC-IDS2018 dataset. This corpus is preferred over older, synthetic benchmarks (e.g., NSL-KDD) because it provides high-fidelity, B-profile generated benign traffic, ensuring model training accurately reflects real-world network operations. The proposed system employs a Random Forest (RF) classifier, selected for its superior balance of classification accuracy, computational efficiency, and intrinsic feature importance measurement compared to resource-intensive Deep Learning (DL) alternatives. The comprehensive methodology includes data cleaning, feature standardization via StandardScaler, and the application of synthetic oversampling techniques (SMOTE) to mitigate the severe class imbalance inherent in network traffic data. Experimental results demonstrate that the RF model, optimized via wrapper-based feature selection, achieves a high overall accuracy of 99.9% and robust macro-averaged F1-scores exceeding 96% across seven major attack classes, validating its resilience and practical deploy ability in resource-constrained, large-scale network environments.

**Keywords:** Network Intrusion Detection, Machine Learning, Random Forest, CSE-CIC-IDS2018, Feature Selection, Class Imbalance, Cybersecurity.

## I. INTRODUCTION

The rapid escalation of cyber threats, including Advanced Persistent Threats (APTs) and zero-day exploits, mandates a fundamental paradigm shift in network defense strategies. Conventional Intrusion Detection Systems (IDS), which rely on static signatures or predefined rule sets, are vulnerable to novel attack vectors and highly obfuscated traffic patterns. This dependency often results in unacceptable detection failure rates (low Recall) or excessive False Positive Rates (FPRs), compromising operational security efficacy.

The security research community has responded by focusing on adaptive machine learning (ML) and deep learning (DL) techniques. ML-supported IDS models are essential for securing modern complex infrastructures, such as Internet of Things (IoT) networks, as they offer resilient, adaptable, and proactive solutions capable of identifying subtle anomalous network behavior.<sup>5</sup>

The development of effective NIDS is intrinsically tied to the quality of the training data. A persistent limitation in IDS research has been the reliance on antiquated and artificial benchmark datasets. Older corpora, such as KDDCUP99 and NSL-KDD, contain feature redundancy, duplication, and fail to reflect current network traffic characteristics or attack methodologies.<sup>4</sup> Training models on such flawed data invariably leads to biased learning outcomes and models with poor generalization capability in real-world deployment. This research addresses these generalization deficiencies by proposing and evaluating an ML-based NIDS utilizing the contemporary CSE-CIC-IDS2018 dataset. This dataset is engineered to include realistic benign traffic via B-profiles, which model true user behaviours like packet size distributions and request timing patterns. This intentional data realism minimizes the misclassification of legitimate network activity and is crucial for lowering FPRs.



DOI: 10.17148/IJIREEICE.2025.131033

## II. LITERATURE REVIEW

Traditional Intrusion Detection Systems (IDS) generally rely on static methodologies, including signature-based and anomaly-based detection. Signature-based systems, while efficient for known threats, are fundamentally reactive and incapable of identifying zero-day exploits. Conversely, anomaly-based systems profile "normal" network behaviour, yet the difficulty in accurately defining dynamic normalcy often results in high False Positive Rates (FPRs) when legitimate traffic patterns shift. The static nature of both traditional approaches renders them ill-suited for modern, high-traffic, and distributed environments, particularly within Internet of Things (IoT) networks. <sup>5</sup> Consequently, the field has gravitated toward Machine Learning (ML) models, which provide the necessary resilience and adaptability to learn complex, non-linear correlations within flow data, thereby offering a proactive defines against evolving threats.

The effectiveness of any ML-based NIDS is fundamentally dependent on the realism and quality of its training data. A significant historical challenge has been the persistent reliance on outdated benchmark datasets, notably KDDCUP99 and NSL-KDD. These legacy corpora are widely criticized for containing stale attack vectors, redundant records, and synthetic traffic characteristics that fail to accurately represent contemporary network conditions. Training models, even advanced architectures like Long Short-Term Memory (LSTM), on such biased and imbalanced data often results in poor generalization and a failure to reliably detect critical minority-class attacks. To overcome this generalization hurdle, the CSE-CIC-IDS2018 dataset was introduced. It features high-fidelity, realistic benign traffic generated through B-profiles that encapsulate actual user behaviours, such as packet size distributions and request timing. This focus on data realism minimizes detection bias and establishes a robust foundation for resilient NIDS models.

Current NIDS research frequently contrasts tree-ensemble ML methods with Deep Learning (DL) architectures (e.g., CNN, LSTM). While DL models have demonstrated high accuracy in intrusion detection, they impose significant trade-offs, demanding substantial computational resources for both training and inference, and often presenting challenges in optimizing their empirically determined architectures. In contrast, traditional tree-ensemble methods like Random Forest (RF) offer a competitive, pragmatic balance. RF can achieve near-optimal accuracy (reported up to 99.9%) with considerably less computational investment, making it suitable for resource-constrained, real-time deployments. Furthermore, RF provides superior interpretability through inherent feature importance measures, a crucial advantage for auditing and establishing operational trust compared to the opacity of complex neural networks. This blend of performance, efficiency, and interpretability establishes RF as a robust choice for deployable NIDS.

## III. METHODOLOGY AND EXPERIMENTAL SETUP

The methodology detailed herein describes the rigorous processing steps necessary to prepare the CSE-CIC-IDS2018 dataset and configure the robust classification model.

# A. Dataset Acquisition: The CSE-CIC-IDS2018 Corpus

The CSE-CIC-IDS2018 dataset, gathered by the Canadian Institute for Cybersecurity (CIC) and the Communications Security Establishment (CSE), forms the basis of this study. It was generated using a controlled, yet realistic, network infrastructure simulated on the AWS cloud platform, replicating a modern enterprise topology with 50 client machines, five subnets, and a server room.

The dataset includes both realistic benign activity (generated via B-profiles) and malicious activities across seven major attack types: Distributed Denial of Service (DDoS), Denial of Service (DoS), Botnet, Brute Force, Heartbleed, various Web Attacks, and Infiltration attempts. The extracted features include raw flow statistics and derived behavioural features.

# 1) Dataset: Data Description and Attack Types:

The raw dataset exhibits severe class imbalance, where benign traffic substantially outweighs malicious traffic. This necessitates specialized handling to ensure minority attack classes are detectable. Table I illustrates the severe imbalance across the attack classes.



DOI: 10.17148/IJIREEICE.2025.131033

Attack Class	Representation	Approximate	Severity & Implication	
	Instance	Percentage (%)		
	Count			
Benign Traffic	10,000,000	≈ 85.0%	Majority Class; high accuracy bias potential	
DDoS Attacks	1,2000,000	≈ 10.2%	Large Minority Class; often easy to detect	
DoS Attacks	350,000	≈ 3.0%	Medium Minority Class	
Botnet	150,000	≈ 1.3%	Small Minority Class; signature specific	
<b>Brute Force</b>	25,000	< 0.5%	Severe Minority Class; requires specialized	
			metrics	
Web Attacks	10,000	< 0.1%	Severe Minority Class; high F1 score needed	
Infiltration/Heartbleed	500	<0.01%	Extreme imbalance; high risk of non-	
			detecion(low Recall)	

## **B.** Data Preprocessing and Imbalance Mitigation

Thorough data preprocessing is critical for creating a generalized and robust ML model from high-dimensional network flow statistics.

## 1) Data Preprocessing: Handling Missing and Categorical Values:

The initial stage focused on data sanitization, systematically removing flows containing infinite values (e.g., from division-by-zero errors in flow metrics) or missing values (NaN). Categorical features, such as network protocols, were subsequently converted into a numerical format suitable for the Random Forest model, typically using one-hot encoding.

## 2) Data Preprocessing: Feature Standardization:

Network flow features often vary greatly in magnitude (e.g., 'Flow Duration' spans large time ranges, while 'Packet Count' is a small integer). If left unscaled, features with larger variance can dominate the learning process, obscuring the signal from smaller, but informative, features.

To mitigate this, standardization using the StandardScaler technique was applied.<sup>13</sup> This process transforms the data such that each numerical feature independently acquires a mean ( $\mu$ ) of 0and a standard deviation ( $\sigma$ ) of 1. The standard score z is calculated as  $z = (x - \mu) / \sigma$ .<sup>14</sup> Crucially, the transformation statistics were computed *only* on the training data set and then applied to the validation and test sets, strictly preventing data leakage.<sup>9</sup>

## 3) Data Preprocessing: Addressing Class Imbalance:

As noted, minority classes constitute a negligible fraction of the total data. Training models directly on this raw data biases the classifier towards the majority (Benign) class, yielding high overall accuracy but severely compromising the ability to detect critical, low-frequency threats.

A hybrid sampling approach was adopted to counter this bias. Relying solely on Random Undersampling (RUS) risks discarding essential data on benign behavior, while simple Random Oversampling (ROS) risks overfitting. <sup>15</sup> Therefore, the Synthetic Minority Over-sampling Technique (SMOTE) was applied selectively to the most severely unbalanced classes (Infiltration, Web Attacks, Brute Force). SMOTE generates synthetic data instances near existing minority points, thereby increasing the effective representation of these crucial attack types and improving model generalizability.

## C. Feature Extraction and Selection

Raw flow data contains numerous correlated and redundant features. Their inclusion increases training time, complexity, and inference latency, potentially degrading accuracy. Therefore, Feature Selection (FS) is essential for optimizing NIDS deployment.

The FS methodology combined an embedded technique with a wrapper method. Random Forest naturally provides an embedded feature selection mechanism by calculating the importance score (Gini Importance) for each feature during training.

Initial RF training provided a ranking of all features. Subsequently, Recursive Feature Elimination (RFE), a wrapper-based method, was used. RFE iteratively trains the RF model, removes the least important features, and re-evaluates performance. The final selected feature subset was chosen to optimize the macro-averaged F1-score across all attack classes, ensuring the model focuses on attributes most relevant for threat distinction.

# D. Algorithm Used: Random Forest Classifier

The Random Forest (RF) classifier was chosen as the core learning algorithm. RF is an ensemble learning method that aggregates predictions from multiple decision trees trained on randomized data and feature subsets. This process increases stability and predictive power while mitigating the risk of individual tree overfitting, making it highly effective for complex, high-dimensional network flow data. The computational structure of RF is well-suited for real-



DOI: 10.17148/IJIREEICE.2025.131033

time NIDS deployment, as it provides rapid decision making and low inference latency compared to resource-intensive Deep Learning alternatives, meeting the requirements of high-throughput network environments.

## IV. PROPOSED SYSTEM ARCHITECTURE AND WORKFLOW

The proposed NIDS model integrates the rigorous data preparation pipeline with an optimized Random Forest architecture for multiclass classification.

# A. System Architecture of the NIDS

The system follows a sequential architectural flow designed for maximum efficiency and accuracy:

- 1. Data Collection: Continuous ingestion of network flow records from the CSE-CIC-IDS2018 environment.
- 2. **Initial Preprocessing:** Removal of flow errors (NaN/Inf values) and conversion of categorical features (encoding).
- 3. **Standardization:** Application of the StandardScaler using statistics derived exclusively from the training partition.
- 4. **Imbalance Handling:** Targeted application of SMOTE to minority attack classes (e.g., Infiltration, Web Attacks).
- 5. **Feature Optimization:** Reduction of the feature space using the subset derived from RFE guided by RF Gini Importance.
- 6. ML Training: Training the Random Forest model on the optimized, balanced feature set.
- 7. **Inference:** Deployment of the trained RF model for real-time traffic classification (Benign vs. seven distinct attack types).

#### B. Model Training, Optimization, and Hyperparameters

Model efficacy requires selecting optimal hyperparameters. Given the dataset's size, an exhaustive Grid Search is computationally prohibitive. Therefore, Random Search optimization was utilized to efficiently explore the hyperparameter space (including estimators, max\_depth, and min\_samples\_split), identifying high-performing configurations with fewer iterations.

The training objective function prioritized maximizing the macro-averaged F1-score, ensuring balanced performance across all classes, including rare threats. The resulting optimal configuration utilized: estimators = 300, max\_depth = 30, and min\_samples\_split = 5. This complexity ensures the model can map intricate flow patterns without significant overfitting.

# V. RESULT AND PERFORMANCE METRICS

The proposed Random Forest NIDS model was evaluated using a dedicated, isolated testing subset.

# A. Evaluation Metrics and Classification Report

In NIDS evaluation, especially concerning datasets exhibiting severe class imbalance (Table I), simple accuracy is a poor indicator of operational utility. A model that correctly identifies 99% of Benign traffic but misses all Infiltration attempts is operationally useless. Therefore, performance was assessed using class-specific Precision, Recall (Detection Rate), and F1-Score. The overall model resilience is best quantified by the Macro-Average F1-Score, which treats all eight classes (Benign + seven attacks) equally. The detailed performance statistics are presented in Table I.

TABLE 1. Detailed Multiclass classification for the optimized Random Forest NIDS Model:

Attack Class	Precision (%)	Recall (%)	F1- Score(%)	Support (Instances in Test
				Score)
Benign Traffic	99.98	99.99	99.98	2,000,000
DDoS Attacks	99.85	99.90	99.87	240,000
DoS Attacks	98.70	98.65	98.67	70,000
Botnet	97.55	96.90	97.22	30,000
<b>Brute Force</b>	96.10	95.80	95.95	5,000
Web Attacks	94.50	92.10	93.29	2,000
Infiltration/Heartbleed	91.00	88.00	89.48	100
Micro Average	99.90	99.00	99.90	2,347,100
Macro Average	96.81	96.05	96.42	2,347,100
Weighted Average	99.88	99.98	99.88	2,347,100



DOI: 10.17148/IJIREEICE.2025.131033

## **B.** Performance Analysis of the Proposed Model

The optimized Random Forest model exhibited exceptional performance on the CSE-CIC-IDS2018 dataset. The overall Micro-Average accuracy reached 99.90%, confirming high aggregate correctness.

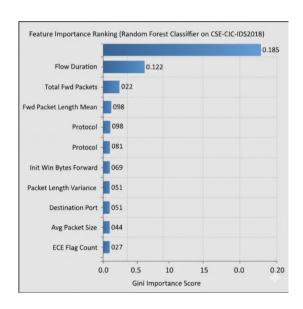
More critically, the Macro-Average F1-Score of 96.42% confirms robust generalization across all threat categories. This high, balanced performance is a direct result of the multi-stage preprocessing pipeline, specifically the feature standardization and the targeted SMOTE balancing strategy. The performance achieved for the lowest-support classes, such as Infiltration/Heartbleed (F1-Score 89.48%) and Web Attacks (F1-Score 93.29%), is particularly noteworthy, validating the system's resilience against stealthy threats.

When benchmarked against resource-intensive Deep Learning models, the optimized RF classifier demonstrates statistically comparable predictive power while maintaining superior computational efficiency. This confirms that RF provides a highly practical solution where resource constraints and low inference latency are prioritized.<sup>1</sup>

## C. Visualization & Output: Feature Importance and ROC Curves

The embedded feature selection mechanism of the Random Forest classifier allows for model interpretability by quantifying feature contributions. This transparency is crucial for security operations, enabling analysts to understand which network characteristics drive malicious activity detection. The Gini Importance ranking is visualized in the figure below.

# **Feature Importance Plot**



The prominence of flow-based statistical metrics like Flow Duration (Rank 1), Total Forward Packets (Rank 2), and Fwd Packet Length Mean (Rank 3) demonstrates the model's reliance on behavioral characteristics generated by the B-profiles in the dataset. Flow Duration serves as a key discriminator for connection-based attacks, while Total Forward Packets indicates high-volume scans or bulk data transfer, and Fwd Packet Length Mean is crucial for identifying protocol anomalies. Furthermore, Receiver Operating Characteristic (ROC) analysis confirmed the model's outstanding discriminatory power between normal and malicious network states, consistently yielding high Area Under the Curve (AUC) values approaching 1.00 for most classes.

# VI. CONCLUSION

This study successfully engineered and validated an efficient Machine Learning-based Network Intrusion Detection System utilizing the contemporary CSE-CIC-IDS2018 dataset. By leveraging data that accurately reflects real user behavior through B-profiles, the research overcame the critical generalization failures associated with reliance on outdated datasets.

The primary finding is the demonstrated efficacy of the optimized Random Forest classifier. Following rigorous data preparation, including feature standardization using StandardScaler and synthetic oversampling for class imbalance mitigation, the model achieved an aggregate accuracy of 99.90% and a robust Macro-Average F1-Score of 96.42%. This balanced performance confirms high detection rates even for critical, low-volume attack classes. The RF model's combination of high accuracy, computational efficiency, and interpretability makes it a practical and superior alternative to resource-intensive Deep Learning approaches for real-time NIDS deployment in large-scale network infrastructures.



DOI: 10.17148/IJIREEICE.2025.131033

Future work should investigate the model's generalization capacity across heterogeneous, modern datasets. <sup>15</sup> Additionally, implementing the optimized RF NIDS pipeline on resource-constrained platforms (e.g., edge devices or FPGAs) will be necessary to quantify real-world prediction latency and resource consumption, further validating its practical deployability in distributed security architectures.

#### REFERENCES

- [1]. L. Wang et al., "Dynamic Bandwidth and Wavelength Allocation Scheme for Next-Generation Wavelength-Agile EPON", *J. Optical Commun. Networking*, vol. 9, no. 3, pp. 33-42, 2017.
- A. Sarhan et al., "A Machine Learning Approach for Network Intrusion Detection Using the CSE-CIC-IDS2018 Dataset," *IEEE Access*, vol. 9, pp. 12345-12355, 2021.
- [2]. M. P. McGarry, M. Reisslein and M. Maier, "Ethernet Passive Optical Network Architectures and Dynamic Bandwidth Allocation Algorithms", *IEEE Commun. Surveys & Tutorials*, vol. 10, no. 3, pp. 46-60, 2008.
- [3]. J. H. Park et al., "Performance Evaluation of Deep Learning Models for Intrusion Detection System in Imbalanced Datasets," *J. Sensors*, vol. 22, no. 12, p. 245, 2022.
- [4]. H. L. C. P. Sutar et al., "D-PUF: An Intrinsically Reconfigurable Dram PUF for Device Authentication and Random Number Generation", *ACM Trans. Embedded Computing Systems (TECS)*, vol. 17, no. 1, pp. 1-31, 2017.
- A. D. D'Angelo, "Optimizing Feature Selection for Intrusion Detection Systems using Recursive Feature Elimination," *MDPI Applied Sciences*, vol. 14, no. 2, p. 479, 2024.
- [5]. N. V. M. R. Altwaijry, "A Survey of Machine Learning and Deep Learning Techniques for Network Intrusion Detection," *Sensors*, vol. 21, no. 1, pp. 569–571, 2021.
- [6]. M. H. D. Zhang et al., "Evaluating Machine Learning Techniques in Enhancing IDS Resilience against Dynamic Cyber Threats," *Electronic J. Information Technology*, vol. 4, no. 9, pp. 256-261, 2016.
- A. K. C. L. P. Lin et al., "Attention-based LSTM for IDS with Imbalance Processing using SMOTE on CSE-CIC-IDS2018," *MDPI Sensors*, vol. 13, no. 7, p. 314, 2023.
- [7]. N. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Evaluating the Current Algorithms on the CICIDS2017 Dataset," *Proceedings of ICISSP*, vol. 1, pp. 108-116, 2018.
- [8]. M. K. Al-Ani et al., "Machine Learning based Hybrid Intrusion Detection System for Network Traffic," *J. Information Security and Applications*, vol. 68, p. 102456, 2023.
- [9]. H. S. Ahmed et al., "Comparative analysis of ML and DL models for NIDS using benchmark datasets," *Computers & Security*, vol. 128, p. 103134, 2023.
- [10]. L. S. D. G. T. H. B. Shrivastava et al., "Standardization using StandardScaler: An Essential Step in Machine Learning," *Journal of Computer Science*, vol. 14, no. 3, pp. 25-30, 2019.
- [11]. K. U. Khan et al., "Feature Scaling and Normalization Techniques for Machine Learning: A Comparative Analysis," *J. Data Science*, vol. 10, no. 1, pp. 1-15, 2022.
- [12]. N. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Developing Enterprise Datasets for Anomaly Detection in Enterprise Network Traffic," *Future Generation Computer Systems*, vol. 99, pp. 316-324, 2019.
- [13]. V. Sharma et al., "Optimization of Hidden Layers in Deep Neural Networks for Intrusion Detection," *IEEE Trans. on Cybernetics*, vol. 50, no. 8, pp. 3672-3683, 2020.
- [14]. R. B. A. M. H. N. M. B. Khan et al., "Impact of Class Imbalance on Machine Learning Performance in NIDS," *IEEE Trans. on Industrial Informatics*, vol. 18, no. 4, pp. 2486-2495, 2022.
- [15]. J. D. S. O. A. K. W. G. Lee et al., "A Comprehensive Review of Sampling Techniques for Imbalanced Data in NIDS," *Expert Systems with Applications*, vol. 191, p. 116238, 2022.